# Sketched Ridge Regression:

## Kernel and Overdetermined Problems

**Alex Gittens**　　many collaborators

RPI　　　　　　　　UC Berkeley

gittea@rpi.edu

**SIAM Annual Meeting 2018 | Portland**

# Kernel Ridge Regression

- Given dataset $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{n}$ and kernel function $\kappa(\mathbf{x}_1, \mathbf{x}_2)$, the problem is to solve

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \|\mathbf{K}\boldsymbol{\alpha} - \mathbf{Y}\|_2^2 + \lambda \boldsymbol{\alpha}^T \mathbf{K}\boldsymbol{\alpha}$$

- Optimal solution:

$$\boldsymbol{\alpha}^\star = (\mathbf{K} + \lambda \mathbf{I}_n)^{-1} \mathbf{Y}$$

- For large $n$ (i.e. $n \approx 10^6$), $\mathbf{K}$ does not even fit in memory

# Iterative Methods

- Since solution doesn't fit in memory, turn to iterative methods
- Classical methods: Conjugate-Gradient, and *Gauss-Siedel*

- We consider randomized block GS (block coordinate descent) for solving positive-definite systems of the form

$$\mathbf{A\alpha}=\mathbf{y}$$

- Given a current iterate

$$(\boldsymbol{\alpha}{\downarrow}k+1\,){\downarrow}\mathrm{J}=(\boldsymbol{\alpha}{\downarrow}k\,){\downarrow}\mathrm{J}-\mathbf{A}{\downarrow}\mathrm{JJ}\,{\uparrow}-1\,(\mathbf{A}\boldsymbol{\alpha}{\downarrow}k-\mathbf{y}){\downarrow}\mathrm{J}$$

# Sampling in Block GS

Two reasonable schemes, given a blocksize $p$:

- **Fixed Partition**: Divide $[n]$ into blocks $J_1, \ldots, J_{n/p}$ blocks ahead of time. During the iterates, randomly choose a block $J_{t_k}$ where $t_k \sim Unif(\{1, \ldots, n/p\})$.

- **Random coordinates**: At each iteration, choose uniformly from the set $\{J \in 2^{[n]} : |J| = p\}$.

# Sampling in Block GS

Fixed partitioning is preferable from a systems perspective (cache locality). Random coordinates suffer from slower memory accesses. Why use random coordinates?

A simple example where the sampling makes a large difference: take

$\mathbf{A}\!\downarrow\!\beta = \mathbf{I} + \beta/n\, \mathbf{1}\mathbf{1}\!\uparrow\!T$

Try GS with $n{=}5000$, $p{=}500$, $\beta{=}1000$.

# Sampling in Block GS

# Convergence of Randomized GS

To understand why the behavior differs, look at the theory of randomized GS

**Theorem.** (Gower and Richtárik, 16)

For all $k \geq 0$,

$$\mathbb{E}\|\boldsymbol{\alpha}_k - \boldsymbol{\alpha}_*\|_{\mathbf{A}} \leq (1-\mu)^{k/2} \|\boldsymbol{\alpha}_0 - \boldsymbol{\alpha}_*\|_{\mathbf{A}},$$

where $\mu = \lambda_{min}(\mathbb{E}[\mathbf{P}_{\mathbf{A}^{1/2}S}])$. Here, the randomized column selection matrix $\boldsymbol{S}$ depends on the choice of sampling scheme.

# Sampling in Block GS

For our example

$$\mathbf{A}_\beta = \mathbf{I} + \frac{\beta}{n}\mathbf{1}\mathbf{1}^T ,$$

$$\mu_{part} = \frac{p}{n} + \beta p$$

$$\mu_{rand} = \mu_{part} + \frac{p-1}{n-1}\frac{\beta p}{n + \beta p}$$

As $\beta \to \infty$, $\mu_{part} \to 1/\beta$ whereas $\mu_{rand} \to p/n$. **This gap is arbitrarily large.**

# Sampling Tradeoffs

- **Systems Perspective**: fixed partition sampling is preferable. Can cache blocks ahead of time, replicate across nodes, etc. *Locality is good for performance.*

- **Optimization perspective**: random coordinates is preferable. Each iteration of GS will make more progress. *Locality is bad for optimization.*

# What about acceleration?

Add a Nesterov momentum step to the iterates.

- Does the same sampling phenomenon occur with acceleration?
- Does this provide the $\sqrt{\mu}$ behavior we expect?

(Assuming the acceleration parameters are carefully chosen)

# Prior State of Theory

The behavior of accelerated **fixed-partition** sampling is understood

**Theorem.**   (Nesterov and Stich, 16)

For all $k \geq 0$, accelerated block GS with fixed-partition sampling satisfies

$$\mathbb{E}\|\boldsymbol{\alpha}_k - \boldsymbol{\alpha}_*\|_{\mathbf{A}} \lesssim (1 - \sqrt{p/n}\,\mu_{part})^{k/2} \|\boldsymbol{\alpha}_0 - \boldsymbol{\alpha}_*\|_{\mathbf{A}},$$

where $\mu_{part} = \lambda_{min}(\mathbb{E}[\mathbf{P}_{\mathbf{A}^{1/2}\mathbf{S}}])$. Here, the randomized column selection matrix $\mathbf{S}$ corresponds to fixed-partition sampling.

Thus fixed-partition sampling loses a factor of $\sqrt{p/n}$ over the ideal Nesterov rate.

# Main Result

**Theorem.**

For all $k \geq 0$, accelerated block GS with any (non-degenerate) sampling scheme satisfies

$$\mathbb{E}\|\boldsymbol{\alpha}_k - \boldsymbol{\alpha}_*\|_A \lesssim (1-\tau)^{k/2} \ \|\boldsymbol{\alpha}_0 - \boldsymbol{\alpha}_*\|_A .$$

Here $\tau = \sqrt{\mu/\nu}$ , where $\mu$ is as before and $\nu$ is a new quantity which behaves roughly like $n/p$.

We prove **this rate is sharp**—there exists a starting point which matches the rate up to constants.

# Corollaries

- For fixed partition sampling, we can show that $\nu = n/p$, recovering Nesterov and Stich's earlier result. Combined with the sharpness of the rate, this proves the $\sqrt{p/n}$ loss over the ideal rate is real for the fixed-partition scheme.

- For random coordinate sampling, we can prove the weaker claim

$$\nu \leq n/p \max_{|J|=p} \frac{\max_{i \in J} \mathbf{A}_{ii}}{\lambda_{min}(\mathbf{A}_{JJ})}$$

If all the size J principal submatrices of **A** are sufficiently well-conditioned, $\nu \approx n/p$.

# Experiment: Accuracy vs Iteration



CIFAR-10 KRR, n=250k, p=10k

# Experiment: Accuracy vs Time



CIFAR-10 KRR, n=250k, p=10k

# Overdetermined Ridge Regression

$$\min_{\mathbf{w}} \; \{ f(\mathbf{w}) = 1/n \, \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \gamma \|\mathbf{w}\|_2^2 \}$$



Applications:

- Basic ML
- IRLS for $\ell_2$-penalized GLMs
- Building block in general optimizers

Two Perspectives:

- (Optimization) Deterministic X, y
- (Statistical) Deterministic X, random y

# Ridge Regression

$$\min_{\mathbf{w}} \left\{ f(\mathbf{w}) = 1/n \, ||\mathbf{Xw} - \mathbf{y}||_2^2 + \gamma ||\mathbf{w}||_2^2 \right\}$$



- **Efficient** and **approximate** solution?
- Use only **part** of the data?

# Ridge Regression

$\min_{\mathbf{w}}\ \{f(\mathbf{w}) = 1/n\ ||\mathbf{Xw} - \mathbf{y}||\downarrow 2 \uparrow 2 + \gamma ||\mathbf{w}||\downarrow 2 \uparrow 2\ \}$



$$\min_{\mathbf{w}}\ \frac{1}{n} \left\| \quad - \quad \right\|_2^2 + \gamma \left\| \quad \right\|_2^2$$

**Matrix Sketching:**

- Random selection
- Random projection

# Approximate Ridge Regression

$\min_{\mathbf{w}} \ \{ f(\mathbf{w}) = 1/n \ ||\mathbf{Xw} - \mathbf{y}||_2^2 + \gamma ||\mathbf{w}||_2^2 \}$

## Optimization Perspective



$\min_{\mathbf{w}} \ \dfrac{1}{n} \left\| \quad \cdot \ - \ \right\|_2^2 + \gamma \left\| \ \right\|_2^2$

s: sketch size

- Sketched solution: $\mathbf{w}^{\uparrow s}$
- $f(\mathbf{w}^{\uparrow s}) \leq (1+\epsilon) \min_{\mathbf{w}} f(\mathbf{w})$

# Approximate Ridge Regression

$\min_{\mathbf{w}} \{ f(\mathbf{w}) = 1/n \|\mathbf{Xw} - \mathbf{y}\|_2^2 + \gamma \|\mathbf{w}\|_2^2 \}$

## Statistical Perspective



- Bias $\|\|_2$ : $\|\mathbf{X}w^\star - \mathbb{E}\mathbf{X}w^s\|$

- Variance $\|\|_2^2$ : $\mathbb{E}\|\mathbf{X}w^s - \mathbb{E}\mathbf{X}w^s\|$

# Related Works on Sketching

Least Squares Regression: $\min_{\mathbf{w}} \|\mathbf{Xw} - \mathbf{y}\|_2^2$

Drineas, Mahoney, and Muthukrishnan. *Sampling algorithms for l2 regression and applications.* SODA, 2006.
Clarkson and Woodruff. *Low rank approximation and regression in input sparsity time.* STOC, 2013.
Raskutti and Mahoney. *A statistical perspective on randomized sketching for ordinary least-squares.* JMLR, 2016.

Ridge Regression: $\min_{\mathbf{w}} 1/n \|\mathbf{Xw} - \mathbf{y}\|_2^2 + \gamma \|\mathbf{w}\|_2^2$

Lu et al. *Faster Ridge Regression via the SRHT.* NIPS, 2013.
Chen et al. *Fast relative-error approximation algorithm for ridge regression.* UAI, 2015.
Avron, Clarkson, Woodruff. *Sharper bounds for Regularized Data Fitting.* Preprint, 2017.
Thanei, Heinze, Meinshausen. *Random projections for large-scale regression.* In Big and Complex Data Analysis, 2017.

# Matrix Sketching



- We consider only efficient sketching procedures
  - Time cost is o$(nds)$ — lower than multiplication.
- Examples:
  - Leverage score sampling: $O(nd\log n)$ time
  - SRHT: $O(nd\log s)$ time

# Sketched Ridge Regression

- Sketched solution:

$$\mathbf{w}\uparrow s = \operatorname*{argmin}_{\mathbf{w}} \{ 1/n \, \|\mathbf{S}\uparrow T \mathbf{X}\mathbf{w} - \mathbf{S}\uparrow T \mathbf{y}\|\downarrow 2 \uparrow 2 + \gamma \|\mathbf{w}\|\downarrow 2 \uparrow 2 \}$$

$$= (\mathbf{X}\uparrow T \mathbf{S}\mathbf{S}\uparrow T \mathbf{X} + n\gamma \mathbf{I}\downarrow d)\uparrow \dagger (\mathbf{X}\uparrow T \mathbf{S}\mathbf{S}\uparrow T \mathbf{y})$$


- Time: $O(sd\uparrow 2) + T\downarrow s$
  - $T\downarrow s$ is the cost of sketching $\mathbf{S}\uparrow T \mathbf{X}$
  - E.g. $T\downarrow s = O(nd \log s)$ for SRHT.
  - E.g. $T\downarrow s = O(nd \log n)$ for leverage score sampling.
- Versus the time for the full RR problem: $O(nd\uparrow 2)$

# Results: Optimization Perspective

# Optimization Perspective

For the sketching methods
- SRHT or leverage sampling with $s=O(\beta d/\epsilon)$,
- uniform sampling with $s=O(\mu \beta d\log d /\epsilon)$,

$f(\mathbf{w}\uparrow s) \leq (1+\epsilon)f(\mathbf{w}\uparrow\star)$ holds w.p. 0.9.

- $\mathbf{X}\in\mathbb{R}\uparrow n\times d$:  the design matrix
- $\gamma$: the regularization parameter
- $\beta=\|\mathbf{X}\|\downarrow 2\uparrow 2 /n\gamma+\|\mathbf{X}\|\downarrow 2\uparrow 2 \in(0, 1]$
- $\mu\in[1,n/d]$:  the row coherence of $\mathbf{X}$

# Optimization Perspective

For the sketching methods
- SRHT or leverage sampling with $s = O(\beta d/\epsilon)$,
- uniform sampling with $s = O(\mu \beta d \log d /\epsilon)$,

$f(\mathbf{w}^{\uparrow s}) \leq (1+\epsilon) f(\mathbf{w}^{\uparrow \star})$ holds w.p. 0.9.

$$\Longrightarrow \quad 1/n \|\mathbf{X}\mathbf{w}^{\uparrow s} - \mathbf{X}\mathbf{w}^{\uparrow \star}\|^{\downarrow 2 \uparrow 2} \leq \epsilon f(\mathbf{w}^{\uparrow \star}).$$

- $\mathbf{X} \in \mathbb{R}^{\uparrow n \times d}$: the design matrix
- $\gamma$: the regularization parameter
- $\beta = \|\mathbf{X}\|^{\downarrow 2 \uparrow 2} /n\gamma + \|\mathbf{X}\|^{\downarrow 2 \uparrow 2} \in (0, 1]$
- $\mu \in [1, n/d]$: the row coherence of $\mathbf{X}$

# Results: Statistical Perspective

# Statistical Model

- $\mathbf{X} \in \mathbb{R}^{n \times d}$ : fixed design matrix

- $\mathbf{w}_0 \in \mathbb{R}^d$ : the *true* and *unknown* model

- $\mathbf{y} = \mathbf{X} \mathbf{w}_0 + \boldsymbol{\delta}$ : observed response vector
  - $\delta_1, \cdots, \delta_n$ are random noise
  - $\mathbb{E}[\boldsymbol{\delta}] = \mathbf{0}$ and $\mathbb{E}[\boldsymbol{\delta}\boldsymbol{\delta}^T] = \xi^2 \mathbf{I}_n$

# Bias-Variance Decomposition

- Risk:  $R(\mathbf{w}) = 1/n \, \mathbb{E} \| \mathbf{X}\mathbf{w} - \mathbf{X}\mathbf{w}_0 \|_2^2$
  - $\mathbb{E}$ is taken w.r.t. the random noise $\boldsymbol{\delta}$.

# Bias-Variance Decomposition

- Risk:    $R(\mathbf{w}) = 1/n \; \mathbb{E} \| \mathbf{Xw} - \mathbf{Xw}_0 \|_2^2$
  - $\mathbb{E}$ is taken w.r.t. the random noise $\boldsymbol{\delta}$.
  - Risk measures prediction error.

# Bias-Variance Decomposition

- Risk:  $R(\mathbf{w}) = 1/n \, \mathbb{E} \lVert \mathbf{Xw} - \mathbf{Xw}_0 \rVert_2^2$

- $R(\mathbf{w}) = \text{bias}^2(\mathbf{w}) + \text{var}(\mathbf{w})$

# Bias-Variance Decomposition

- Risk: $R(\mathbf{w}) = 1/n\, \mathbb{E}\|\mathbf{Xw} - \mathbf{Xw}_0\|_2^2$

- $R(\mathbf{w}) = \text{bias}^2(\mathbf{w}) + \text{var}(\mathbf{w})$

  - $\text{bias}(\mathbf{w}^\star) = \gamma\sqrt{n}\,\|(\mathbf{\Sigma}^2 + n\gamma\mathbf{I}_d)^{-1}\mathbf{\Sigma V}^T \mathbf{w}_0\|_2,$
  - $\text{var}(\mathbf{w}^\star) = \xi^2/n\,\|(\mathbf{I}_d + n\gamma\mathbf{\Sigma}^{-2})^{-1}\|_2^2,$

  - $\text{bias}(\mathbf{w}_s) = \gamma\sqrt{n}\,\|(\mathbf{\Sigma U}^T \mathbf{SS}^T \mathbf{U\Sigma} + n\gamma\mathbf{I}_d)^\dagger\, \mathbf{\Sigma V}^T \mathbf{w}_0\|_2,$

  - $\text{var}(\mathbf{w}_s) = \xi^2/n\,\|(\mathbf{U}^T \mathbf{SS}^T \mathbf{U} + n\gamma\mathbf{\Sigma}^{-2})^\dagger\, \mathbf{U}^T \mathbf{SS}^T\|_2^2,$

  - Here $\mathbf{X} = \mathbf{U\Sigma V}^T$ is the SVD.

# Statistical Perspective

For the sketching methods

- SRHT or leverage sampling with $s = O(d/\epsilon^2)$;
- uniform sampling with $s = O(\mu\, d \log d /\epsilon^2)$,

<div style="border:1px solid purple">

- $\mathbf{X} \in \mathbb{R}^{n \times d}$ : the design matrix
- $\mu \in [1, n/d]$: the row coherence of $\mathbf{X}$

</div>

the following hold w.p. 0.9:

$$1 - \epsilon \leq \text{bias}(\mathbf{w}^s)/\text{bias}(\mathbf{w}^\star) \leq 1 + \epsilon,$$

Good!

$$(1 - \epsilon)n/s \leq \text{var}(\mathbf{w}^s)/\text{var}(\mathbf{w}^\star) \leq (1 + \epsilon)n/s.$$

Bad! Because $n \gg s$.

# Statistical Perspective

For the sketching methods

- SRHT or leverage sampling with $s = O\left(d / \epsilon^2\right)$;
- uniform sampling with $s = O\left(\mu\, d \log d / \epsilon^2\right)$,

> $\mathbf{X} \in \mathbb{R}^{n \times d}$ : the design matrix
> $\mu \in [1, n/d]$ : the row coherence of $\mathbf{X}$

the following hold w.p. 0.9:

$$1 - \epsilon \leq \mathrm{bias}(\mathbf{w}^s) / \mathrm{bias}(\mathbf{w}^\star) \leq 1 + \epsilon,$$

$$(1 - \epsilon)\, n/s \leq \mathrm{var}(\mathbf{w}^s) / \mathrm{var}(\mathbf{w}^\star) \leq (1 + \epsilon)\, n/s.$$

> If **y** is noisy
> ⟹ variance dominates bias
> ⟹ $R(\mathbf{w}^s) \gg R(\mathbf{w}^\star)$.

# Consequence for selection of regularization



Legend: Uniform Sampling, Leverage Sampling, Shrinked Lev. Sampling, Gaussian Projection, SRFT, Count Sketch, Optimal Solution

$W^c$

$\text{Bias}^2$

Var

$\text{Risk} = \text{Bias}^2 + \text{Var}$

the min risk of the classical sketch

the min risk of the optimal solution

optimal $\gamma$ for the optimal solution

optimal $\gamma$ for the classical sketch

# Model Averaging to Reduce Variance

# Model Averaging

- Independently draw $\mathbf{S}_1, \cdots, \mathbf{S}_g$.

- Compute the sketched solutions $\mathbf{w}_1^s, \cdots, \mathbf{w}_g^s$.

- Model averaging: $\mathbf{w}^s = 1/g \sum_{i=1}^{g} \mathbf{w}_i^s$.

# Connection to Bagging

- Bagging (bootstrap aggregation) was proposed by Breiman in 1996 for reducing the variance of the decision tree.

- Bagging originates in decision tree methods, but it can be used with many machine learning models.

- For ridge regression, uniform sampling with model averaging is exactly bagging.

- Our approach is not limited to uniform sampling. Random projections and non-uniform sampling outperform uniform sampling.

# Optimization Perspective

- For sufficiently large $s$,

  $$f(\mathbf{w}_{1}^{s}) - f(\mathbf{w}^{\star}) / f(\mathbf{w}^{\star}) \leq \epsilon \quad \text{holds w.h.p.}$$

  **Without** model averaging

- Using the **same** sketching distribution and $s$,

  $$f(\mathbf{w}^{s}) - f(\mathbf{w}^{\star}) / f(\mathbf{w}^{\star}) \leq \epsilon/g + \epsilon^{2} \quad \text{holds w.h.p.}$$

  **With** model averaging

# Statistical Perspective

- For sufficiently large $s$, the following hold w.h.p.:

$$\text{bias}(\mathbf{w}_s)/\text{bias}(\mathbf{w}_\star) \leq 1+\epsilon \quad \text{and} \quad \text{var}(\mathbf{w}_s)/\text{var}(\mathbf{w}_\star) \leq n/s \ (1+\epsilon).$$

- Using the **same** sketching distribution and $s$, the following hold w.h.p.:

$$\text{bias}(\mathbf{w}_s)/\text{bias}(\mathbf{w}_\star) \leq 1+\epsilon \quad \text{and} \quad \text{var}(\mathbf{w}_s)/\text{var}(\mathbf{w}_\star) \lesssim n/s \ (1/\sqrt{g} +\epsilon)^2$$

# Empirical variance reduction

- If $s$ is large compared to $d$ and $g$ is larger than $n/s$, then $\text{var}(\mathbf{w}^s) < \text{var}(\mathbf{w}^\star)$.



Experiments on synthetic data.

- $n=10^5$, $d=500$, $\kappa(X^T X)=10^{12}$.
- Sketch size is $s=5000=n/20$.
- Regularization parameter $\gamma=10^{-6}$.
- As $g$ exceeds $n/s=20$, $\text{var}(\mathbf{w}^s)$ can be smaller than $\text{var}(\mathbf{w}^\star)$.

# Thank You!

*Breaking Locality Accelerates Block Gauss-Seidel*. Tu, G., et al. ICML 2017
https://arxiv.org/abs/1701.03863


S. Wang, G., and M. W. Mahoney. "*Sketched Ridge Regression: Optimization Perspective, Statistical Perspective, and Model Averaging*". ICML, 2017.
https://arxiv.org/abs/1702.04837